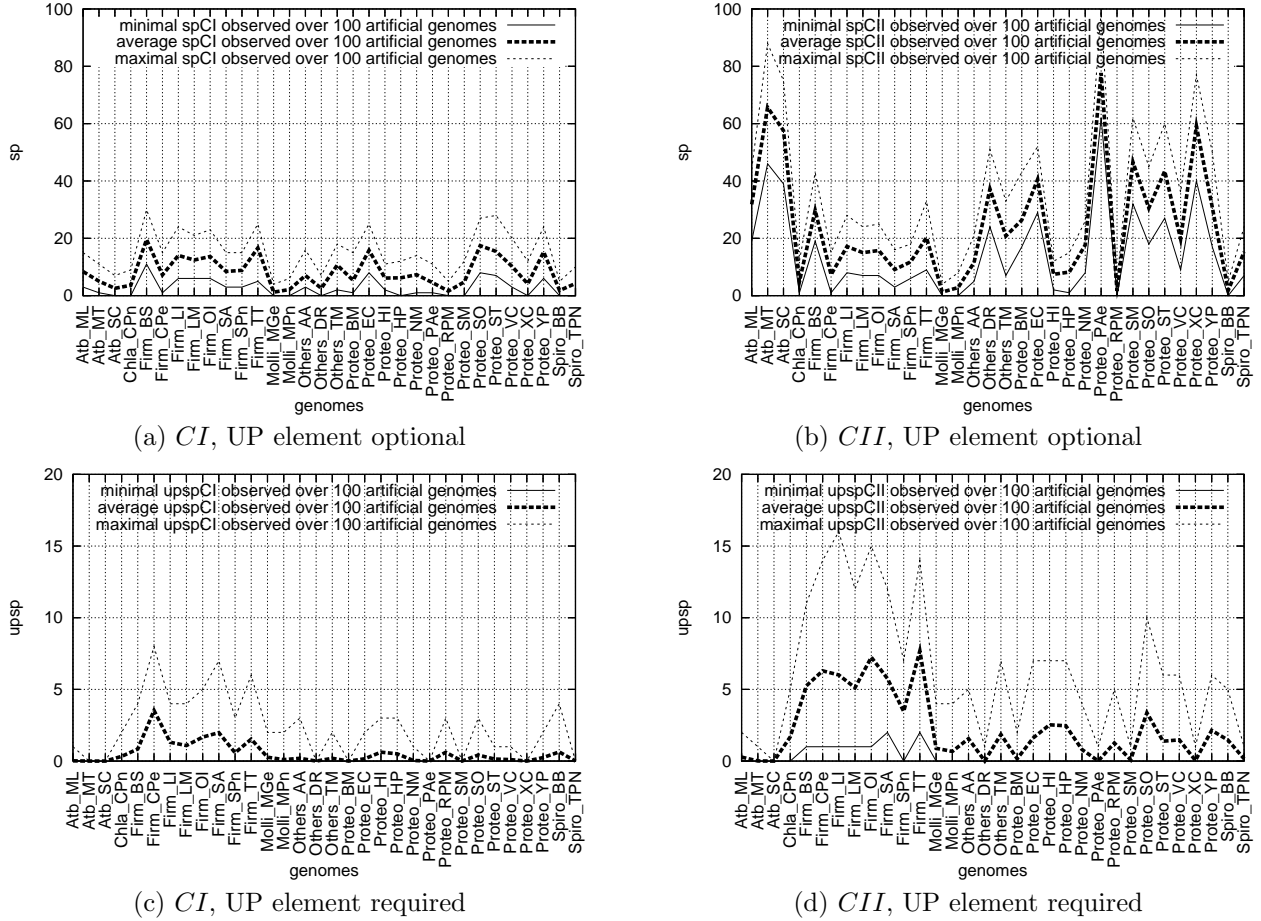


## Appendix 5.

### Comparing observations in bacterial genomes with expectations in randomly generated genomes

#### Magnification of the results obtained for randomly generated genomes



**Figure 5.1** Minimal, average and maximal values observed over 100 randomly generated genomes, for *sp* and *upsp* respectively; 32 bacterial genomes are considered (see Figure 1 for genome nomenclature). For each such bacterial genome, 100 artificial genomes are generated at random, which have same proportions of A, C, T, G nucleotides and same total number of genes encoding proteins as the bacterial genome.  $sp_{CI}$  denotes the number of genes with a putative **Strong Promoter** identified under constraint set *CI*.  $sp_{CII}$  is defined similarly for constraint set *CII* (see Section "Systems and methods", Subsection "Genome analysis upon request" for the definition of *CI* and *CII* constraints).  $upsp_{CI}$  denotes the number of genes with an **UP** element in their putative **Strong Promoter**, and identified under constraint set *CI*.  $upsp_{CII}$  is defined similarly.

#### Description of statistical significance through Z-scores

Table 5.1 compares the Z-scores obtained under all four conditions. When Z-scores could not be calculated, we noticed that the *sp* value (respectively the *upsp* value) observed for the bacterial genome and the corresponding value expected for the average artificial genome are both equal to 0 (exceptionally, such previous values are equal to 1 and 0 respectively).

	<i>CI</i>				<i>CII</i>				<i>CI</i>				<i>CII</i>			
	UP element optional				UP element required											
	min	max	av	std	min	max	av	std	min	max	av	std	min	max	av	std
8 <i>Firmicutes</i>	81.3	308.5	193.0	66.1	92.9	324.6	216.9	66.8	74.9	291.7	155.9	64.3	69.4	311.9	188.3	66.9
13 <i>Proteobacteria</i>	1.0	32.4	16.0	9.5	3.4	53.6	23.5	15.2	2.0	15.6	11.0	7.9	4.4	66.6	16.8	17.0
									(1)							
26 species with large genomes	1.0	308.5	74.9	89.9	4.9	324.6	91.4	97.0	0.15	291.7	76.9	81.6	0.2	311.9	75.9	88.8
									(2)				(3)			

**Table 5.1** Evaluation of the statistical significance of  $\sigma 70$  promoter frequencies: comparison between different bacterial groups. For details about Z-scores, see text, Subsection "Comparing observations in bacterial genomes with expectations in randomly generated genomes"; av: average, std: standard deviation. (1)(2)(3): Z-scores were calculable respectively for 10, 19 and 25 genomes.

genome name	abbreviation	<i>CI</i>				<i>CII</i>				<i>CI</i>				<i>CII</i>			
		UP element optional				UP element required				UP element required				UP element required			
genomes for which Z-score value is above threshold																	
		<b>7</b>	<b>15</b>	<b>80</b>	<b>7</b>	<b>15</b>	<b>80</b>	<b>7</b>	<b>15</b>	<b>80</b>	<b>7</b>	<b>15</b>	<b>80</b>	<b>7</b>	<b>15</b>	<b>80</b>	
<i>Mycobacterium leprae</i> tn	Atb_ML	ML	ML		ML	ML											
<i>Mycobacterium tuberculosis</i> h37rv	Atb_MT	MT			MT	MT		-	-	-	MT	MT					
<i>Streptomyces coelicolor</i> a3 (2)	Atb_SC	SC			SC	SC		-	-	-	-	-	-				
<i>Aquifex aeolicus</i> vf5	Others_AA	AA	AA	AA	AA	AA	AA	AA	AA	AA	AA	AA	AA	AA	AA	AA	
<i>Deinococcus radiodurans</i> r1	Others_DR	DR			DR	DR		-	-	-	DR	AA					
<i>Thermotoga maritima</i>	Others_TM	TM	TM	TM	TM	TM	TM	TM	TM	TM	TM	TM	TM	TM	TM	TM	
<i>Brucella melitensis</i> 16m	Proteo_BM	BM			BM			-	-	-	BM						
<i>Escherichia coli</i> k12	Proteo_EC	EC	EC		EC	EC		EC			EC	EC					
<i>Haemophilus influenza</i> rd kw20	Proteo_HI	HI			HI			HI			HI						
<i>Helicobacter pylori</i> j99	Proteo_HP	HP			HP			HP			HP						
<i>Neisseria meningitidis</i> mc58	Proteo_NM	NM	NM		NM	NM		NM	NM		NM						
<i>Pseudomonas aeruginosa</i> pa01	Proteo_PAe	PAe	PAe		PAe	PAe		-	-	-	PAe	PAe					
<i>Sinorhizobium meliloti</i> 1021	Proteo_SM	SM	SM		SM	SM		-	-	-	SM	SM					
<i>Shewanella oneidensis</i> mr1	Proteo_SO	SO	SO		SO	SO		SO			SO	SO					
<i>Salmonella typhimurium</i> lt2	Proteo_ST	ST	ST		ST	ST		ST	ST		ST	ST					
<i>Vibrio cholerae</i> n16961	Proteo_VC				VC												
<i>Xanthomonas campestris</i> atcc 33913	Proteo_XC							-	-	-							
<i>Yersinia pestis</i>	Proteo_YP	YP			YP	YP					YP						

**Table 5.2** Evaluation of the statistical significance of  $\sigma 70$  promoter frequencies for 18 bacterial species selected as non *Firmicutes* species with large genomes. For definition of *CI* and *CII* constraints, see Section "Systems and methods", Subsection "Genome analysis upon request". Three significance thresholds are considered (7, 15 and 80). – means that the Z-score is not calculable. In a given column, the mention of a species points out statistical significance for the corresponding threshold.

		Z-scores				
	UP element presence	<i>E. coli</i>	min	max	average	standard deviation
<i>CI</i>	optional	<b>21.7</b>	1.0	308.5	74.9	89.9
<i>CII</i>	"	<b>38.2</b>	4.9	324.6	91.4	97.0
<i>CI</i>	required	<b>7.3</b>	0.15	291.7	76.9	81.6
<i>CII</i>	"	<b>15.3</b>	0.2	311.9	75.9	88.8

**Table 5.3** Evaluation of the statistical significance of  $\sigma 70$  promoter frequencies: comparison of *E. coli* with respect to 26 species with large genomes.

		Z-score threshold			number of genomes with calculable Z-scores
	UP element presence	7	15	80	
<i>CI</i>	optional	24	15	10	26
<i>CII</i>	"	25	21	10	26
<i>CI</i>	required	16	12	7	19
<i>CII</i>	"	22	16	8	25

**Table 5.4** Numbers of large genomes (among 26) for which the total number of putative strong promoters identified is shown to be significantly different from that of the corresponding "average" randomly generated genome.

## Is there a genome size bias? Comparison of *Firmicutes* genomes with similarly AT-rich *Proteobacteria* genomes

Hereafter, we complete our analysis, checking that the specificity observed for *Firmicutes* is not due to genome size bias. For this purpose, we compare two *Proteobacteria* genomes (*H. influenza*, *H. pylori*) and four *Firmicutes* genomes (*L. innocua*, *L. monocytogenes*, *S. pneumoniae*, *T. tengcongensis*). All six (high) AT-richnesses range in the narrow interval [60.8%, 62.6%]. Successively considering the number of promoters observed in each *Proteobacteria* genome as a reference, we calculate a corrected frequency for each *Firmicute* genome, applying a correction based on proportionality relative to genome sizes. Then we compare observed values *versus* corrected values under each of the four conditions studied. Not only do we implement such corrections for bacterial genomes, we also process the six average random genomes in a similar way (see Tables 5.5 and 5.6). When considering the average genomes generated at random corresponding to *Firmicutes*, we show that the order of magnitude is identical for expected values and corrected expected values. In contrast, when dealing with bacterial genomes, the size bias correction does not smooth out the difference between *Firmicutes* and *Proteobacteria*.

<i>Proteobacteria</i>			<i>Firmicutes</i>							
	HI	HP	LI		LM		SPn		TT	
AT-content	61.9%	60.8%	62.6%		62.0%		60.28%		62.4%	
g	<b>1673</b>	<b>1478</b>	<b>2962</b>	<i>HI - HP</i>	<b>2837</b>	<i>HI - HP</i>	<b>1861</b>	<i>HI - HP</i>	<b>2588</b>	<i>HI - HP</i>
	obs	obs	obs	corr	obs	corr	obs	corr	obs	corr
<i>sp<sub>CI</sub></i>	31	31	713	<i>54-62(*)</i>	707	<i>52-59</i>	213	<i>34-39</i>	581	<i>47-54</i>
<i>sp<sub>CII</sub></i>	37	34	946	<i>65-68</i>	926	<i>62-65</i>	285	<i>41-42</i>	715	<i>57-59</i>
<i>upsp<sub>CI</sub></i>	8	7	145	<i>14-14</i>	147	<i>13-13</i>	59	<i>8-8</i>	150	<i>12-12</i>
<i>upsp<sub>CII</sub></i>	20	17	501	<i>35-34</i>	488	<i>33-32</i>	120	<i>22-21</i>	407	<i>30-29</i>

**Table 5.5** Comparison of the observed numbers of putative strong promoters between two *Proteobacteria* and four *Firmicutes* genomes characterized by close (high) AT-contents (range [60.2%, 62.4%]), under conditions *CI* and *CII*, and with or without UP element required. *g*: total number of genes encoding proteins in the genome considered; HI: *Haemophilus influenzae*, HP: *Helicobacter pylori*; LI: *Listeria innocua*, LM: *Listeria monocytogenes*, SPn: *Streptococcus pneumoniae*, TT: *Thermoanaerobacter tengcongensis*; obs: observed values, corr: corrected values. (\*) With HI then HP taken as a reference, these columns in italics yield the corrected values *sp<sub>CI</sub>*, ..., *upsp<sub>CII</sub>* on the basis of proportionality to total gene number; the two corrected values in italics have to be compared with the value on their left. Example: the reference being HI, corrected *sp<sub>CI,corr</sub>* for LI is  $sp_{CI,corr}(LI) = \frac{sp_{CI,obs}(HI) \times 2962}{1673} = 54$ , with  $sp_{CI,obs}(HI) = 31$ ;  $sp_{CI,corr}(LI)$  has to be compared with the value of 713 observed for LI.

<i>Proteobacteria</i>			<i>Firmicutes</i>							
	HI	HP	LI		LM		SPn		TT	
AT-content	61.9%	60.8%	62.6%		62.0%		60.28%		62.4%	
g	<b>1673</b>	<b>1478</b>	<b>2962</b>	<i>HI - HP</i>	<b>2837</b>	<i>HI - HP</i>	<b>1861</b>	<i>HI - HP</i>	<b>2588</b>	<i>HI - HP</i>
	exp	exp	exp	corr	exp	corr	exp	corr	exp	corr
<i>sp<sub>CI</sub></i>	6	6	14	<i>10-12(*)</i>	12	<i>10-11</i>	8	<i>6-7</i>	16	<i>9-10</i>
<i>sp<sub>CII</sub></i>	7	8	17	<i>12-16</i>	14	<i>11-15</i>	11	<i>7-10</i>	20	<i>10-12</i>
<i>upsp<sub>CI</sub></i>	0	0	1	<i>0-0</i>	1	<i>0-0</i>	0	<i>0-0</i>	1	<i>0-0</i>
<i>upsp<sub>CII</sub></i>	2	2	6	<i>3-4</i>	5	<i>3-3</i>	3	<i>2-2</i>	7	<i>3-3</i>

**Table 5.6** Comparison of the expected numbers of putative strong promoters between six average genomes generated at random and characterized by the same AT-contents as two *Proteobacteria* and four *Firmicutes* genomes (range [60.2%, 62.4%]), under conditions *CI* and *CII*, and with or without UP element required. See Table 5.5 caption for explanations. Example: the reference being HI, corrected *sp<sub>CI,corr</sub>* for LI is  $sp_{CI,corr}(LI) = \frac{sp_{CI,exp}(HI) \times 2962}{1673} = 10$ , with  $sp_{CI,exp}(HI) = 6$ ;  $sp_{CI,corr}(LI)$  has to be compared with the value of 14 expected for the average artificial genome having the same AT-richness as the LI genome.